# Post-processing of enrichment data
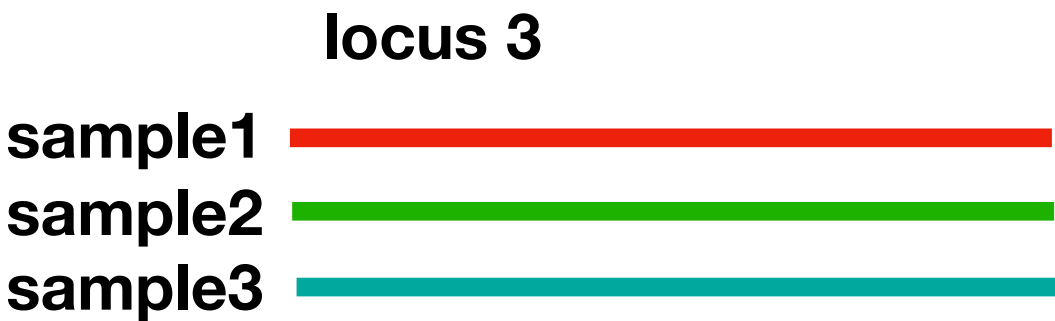
**Reporter: Hao Yuan**

# Output of assembling pipeline

**locus 1**

sample1 ━━━━━━━━━━━━━━━
sample2 ━━━━━━━━━━━━━━━
sample3 ━━━━━━━━━━━━━━━

**locus 2**

sample1 ━━━━━━━━━━━━━━━
sample2 ━━━━━━━━━━━━━━━
sample3 ━━━━━━━━━━━━━━━

**locus 3**

sample1 ━━━━━━━━━━━━━━━
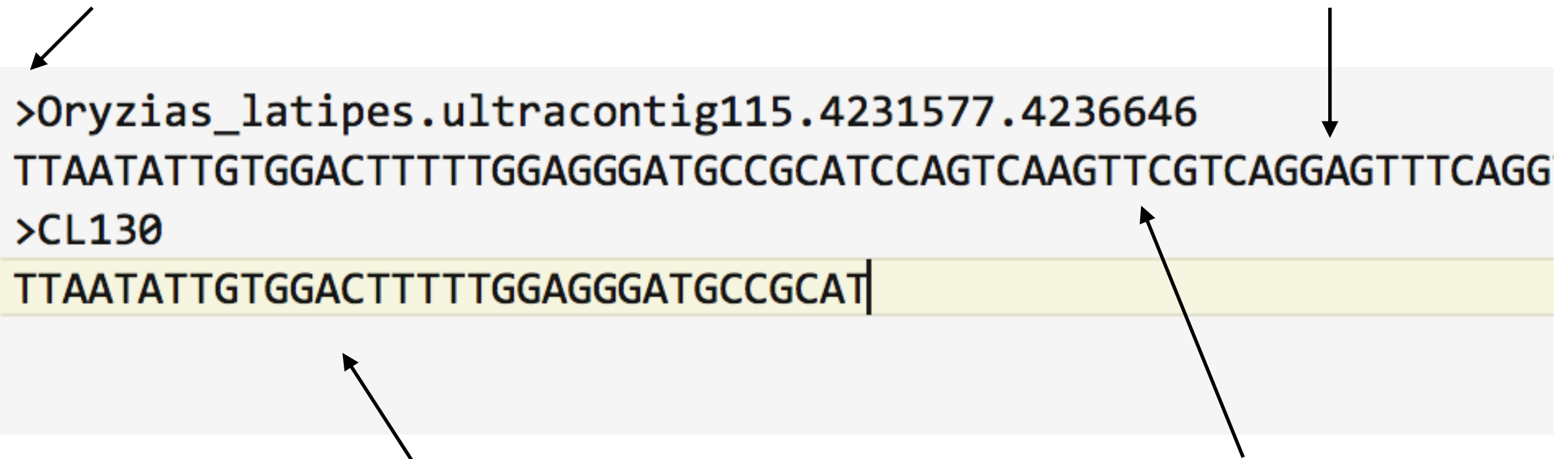sample2 ━━━━━━━━━━━━━━━
sample3 ━━━━━━━━━━━━━━━

**Sequences targeting the same loci**

# Output until here looks like this

A ">" before the sequence name

After name is sequence

```
>Oryzias_latipes.ultracontig115.4231577.4236646
TTAATATTGTGGACTTTTTGGAGGGATGCCGCATCCAGTCAAGTTCGTCAGGAGTTTCAGG
>CL130
TTAATATTGTGGACTTTTTGGAGGGATGCCGCAT
```

First one is the sequence of reference

Following is the sequence of enriched sample

This called **fasta** format. File suffix is "**fa**", "**fas**" or "**fasta**"

**Assembled contigs are coding sequences**

**Start from the
first codon**

**End at  from the
third codon**

**TAT**<span style="color:red">**AAC**</span><span style="color:blue">**CTG**</span>

**Length of nucleotide can
be exactly divided by 3**

**Y** <span style="color:red">**N**</span> <span style="color:blue">**L**</span>

**No stop codon in amino acid sequences**

**Output will be placed under "assemble_result" including 3 folders:**

**1) nf: folder containing full coding nucleotide sequences**

**2) f: folder containing coding sequences with flankings**

**3) p: folder containing amino acid sequences**

# Further processing

Data manipulation

Multiple sequences alignment

Filter poorly aligned regions

Filter for other purpose

Detect cross-sample contamination

Summary statistics

# Data manipulation

**Remove poorly enriched taxa**

**enriched_gene.txt**

```
total    4435

Sample  Num. of enriched genes  Percentage of enriched genes(%)
sample1 3262     73.6
sample2 3253     77.3
sample3 3356     75.7
sample4 3516     79.3
sample5 410 9.24
```

**"—deselected_taxa" option of pick_taxa.pl**

**Get loci with certain data completeness level**

**"—min_seq" option of pick_taxa.pl**

$$10*0.8 = 8$$

number of
total sample

completeness
level

# Data manipulation

## Extract loci from existing genomes



Organism 1 | Organism 2

one of locus
in reference

orthologous
locus in targeted
genome

use merge_loci.pl to add extracted
sequences to enriched dataset

get_orthologues.pl
merge_loci.pl

# Multiple sequences alignment

Reads → Contigs

specify "—non_codon_aln" option of mafft_aln.pl if align sequences with flanks



→ Tree reconstruction

**mafft_aln.pl**

**Multiple sequence alignment (MSA)**

**Align AA/DNA with common ancestry**

# Why need to filter resulting alignment

**Short or too conserved alignment** → **Few phylogenetic signal**

**Low node support**

**Poor alignment** → **Conflict phylogenetic signal**

**Discordant tree using different method**

**Mis-alignment**

**Missing data**

```
ATGCGTACGTT
ATGCGTATT——
```

**Paralogs**

```
ATGCGTACGTT
ATGCGTA—–TT
```

**Saturated mutation**

# Why need to filter resulting alignment

**Short or too conserved alignment**

**Few phylogenetic signal**

**Low node support**

**Poor alignment** → **Conflict phylogenetic signal**

**Discordant tree using different method**

**Mis-alignment**

**Missing data**

**Paralogs**

**Saturated mutation**

# Filter poorly aligned regions

## Poorly aligned coding regions

**Remove sequences having long insertion or deletion to reference**

↓

**Remove sequences distant from reference**

↓

**Remove sequences with low coverage**

↓

**Realign**

filter.pl

# Filter poorly aligned regions

## Sequences distant from reference

| | | | |
|---|---|---|---|
| ref | CAGGACG | TTTACTGAAGTTTTCAGGAGA | GCTTGAAGGTGTTTCAAGAGAAGACT |
| sample1 | CTGGATG | TCTGTTGAAGTTTTCTGGGGT | GCTTGAAGACGTTTCAAGAGAGGACT |
| sample2 | CTGGATG | TCTGTTGAAGTTTTCTGGGGA | GCTTGAAGACGTTTCAAGAGAGGACT |
| sample3 | CTGGATG | TCTGTTGAAGTTTTCTGGGGA | GCTTGAAGATGTTTCAAGAGAGGACT |
| sample4 | CTGGATG | TCTGTTGAAGTTTTCTGGGGA | GCTTGAAGATGTTTCAAGAGAGGACT |
| sample5 | CTGGATG | TCTGTTGAAGTTTTCTGGGGA | GCTTGAAGATGTTTCAAGAGAGGACT |
| sample6 | CTGGATG | TCTGTTGAAGTTTTCTGGGGA | GCTTGAAGATGTTTCAAGAGAGGACT |

**50 bp, 25 bp per step**

# Filter poorly aligned regions

## Poorly aligned coding regions

Remove sequences having long insertion
or deletion to reference

↓

Remove sequences distant from reference

↓

Remove sequences with low coverage

↓

Realign

# Filter poorly aligned regions

## Poorly aligned flanking regions

**Remove unevenly enriched sequences from ends**

$\downarrow$

**Remove sequences having long and unique insertion**

$\downarrow$

**Remove too variable sequences**

$\downarrow$

**Remove short flanks and flanking sequences with low coverage**

$\downarrow$

**Realign**

**flank_filter.pl**

# Filter poorly aligned regions

**Poorly aligned flanking regions**

**Remove unevenly enriched sequences from ends**

↓

**Remove sequences having long and unique insertion**

↓

**Remove too variable sequences**

↓

**Remove short flanks and flanking sequences with low coverage**

↓

**Realign**

# Filter poorly aligned regions

## Remove sequences having long and unique insertion



**chimeric assembly**

**>= 10bp**

# Filter poorly aligned regions

**Poorly aligned flanking regions**

**Remove unevenly enriched sequences from ends**

↓

**Remove sequences having long and unique insertion**

↓

**Remove too variable sequences**

↓

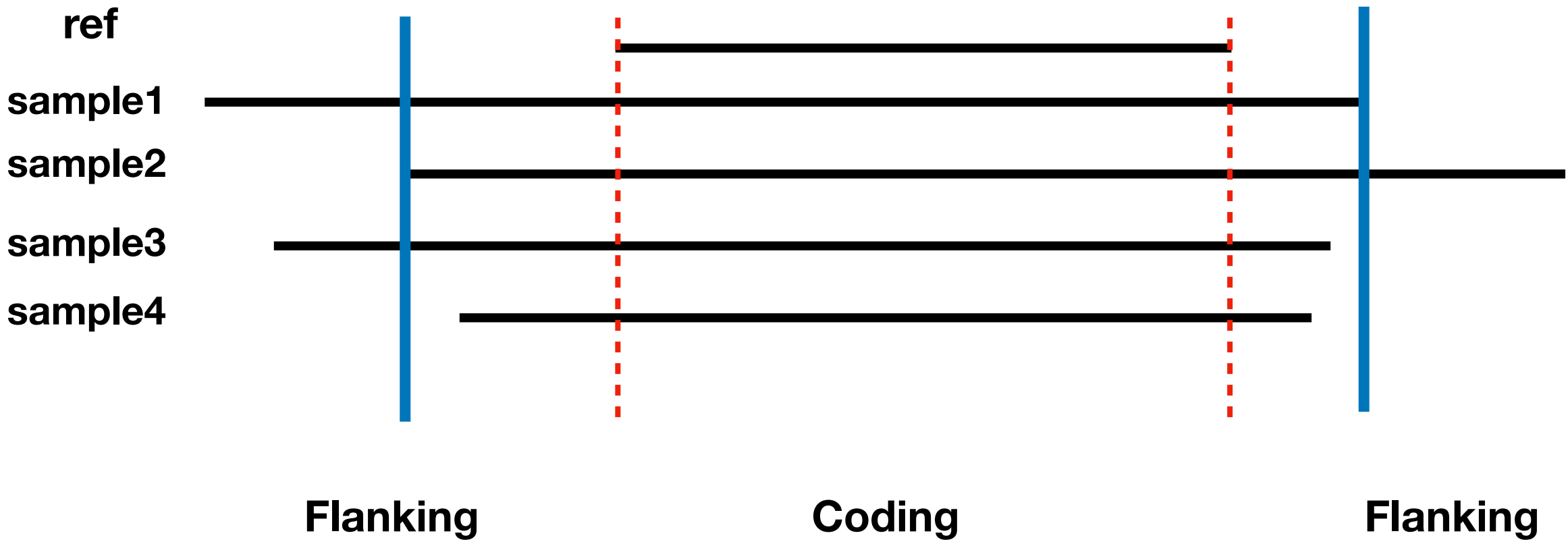**Remove short flanks and flanking sequences with low coverage**

↓

**Realign**

# Filter poorly aligned regions

## Remove too variable sequences

**No reference in flanks**

# Filter poorly aligned regions

## Remove too variable sequences



**consensus reference**

ref

sample1

sample2

sample3

sample4

**Remove sequences distant from consensus reference**

# Filter poorly aligned regions

**Remove too variable sequences**

**consensus reference**

ref

sample1

sample2

sample3

sample4

**2 or more patterns of sequences existing in flanks**

# Filter poorly aligned regions

## Remove too variable sequences



**consensus reference**

ref

sample1

0.45

sample2    0.43

sample3

0.40

sample4

**Find pair of distant sequences, then compute their distance to the rest of the sequences**

# Filter poorly aligned regions

**Remove too variable sequences**



**consensus reference**

ref

sample1

sample2

sample3

sample4

**Remove sequences distant from rest of the sequences**

# Filter poorly aligned regions

## Remove too variable sequences



consensus reference

**Remove sequences distant from rest of the sequences**

# Filter poorly aligned regions

**Remove too variable sequences**

**consensus reference**



**Remove sequences distant from rest of the sequences**

# Filter poorly aligned regions

## Poorly aligned flanking regions

**Remove unevenly enriched sequences from ends**

↓

**Remove sequences having long and unique insertion**

↓

Remove flanks distant from rest of sequences

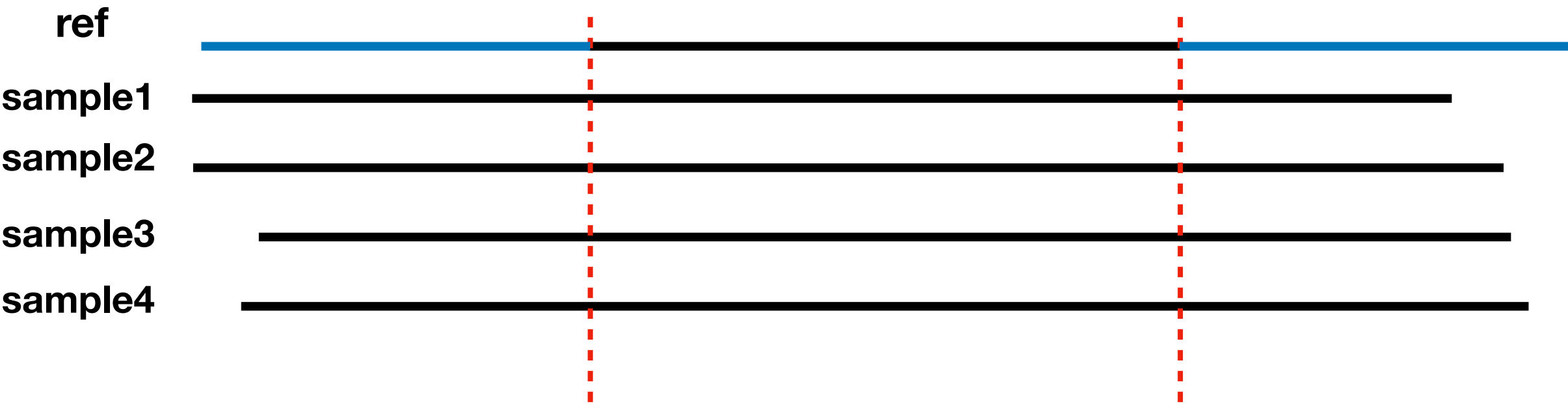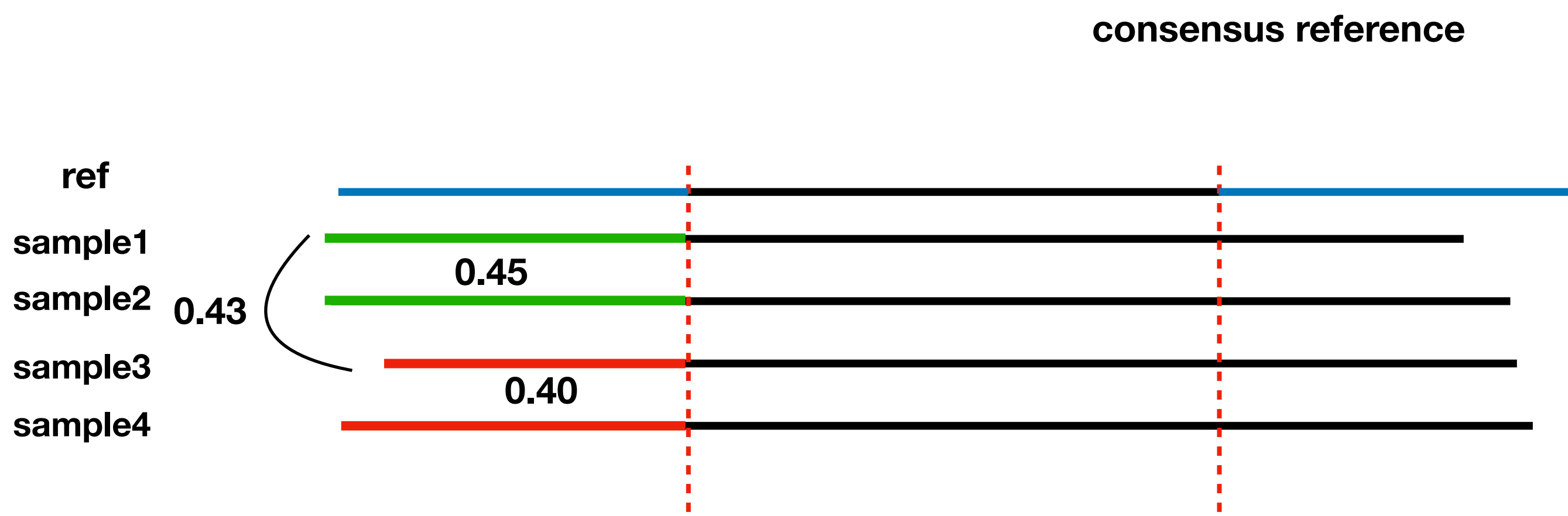**Remove too variable sequences**

Remove sequences distant from consensus reference

↓

**Remove short flanks and flanking sequences with low coverage**

↓

**Realign**

# Filter poorly aligned regions

**Poorly aligned flanking regions**

**Usage:**

```
$ perl flank_filter.pl \
--flank assemble_result/f \            ←————— sequences with flanks
--nonflank_filtered filtered_nf \      ←————— filtered coding sequences
--flank_filtered Oreochromis_niloticus \  ←— filtered sequences with flanks
--ref_taxa Oreochromis_niloticus       ←————— name of reference
```

# Filter for other purpose

**Filter out loci with pre-defined monophyletic group**

**Filter out loci follow the molecular clock hypotheses**

# Filter for other purpose

**Filter out loci with pre-defined monophyletic group**

```
sample1 sample2 sample3 sample4 sample5 sample6
sample7 sample9
sample8 sample10 sample3
sample11 sample12
```

**Define monophyletic group in a txt file, one group a line**

**At least two sample in a group**

# Filter for other purpose

## SH and AU test

ATCGTAGGGCTGGCTAGTCGTAGCTA
ATCGTAGGGCTGGCGAGTCGT–GCTA
ATCGTAGGACTGGCTAGTCGTAGCTA
ATCGTACGGCTGGCTAGTCGTAGCTA
ATCGTAGGGCTGGCTAGTCGTAGCTA

**ML tree**

**Conserved** **Less conserved**

**SH or AU test**

**whether p > 0.05**

**Monophyletic group constrained ML tree**

**No**

**Yes**

**use construct_tree.pl to build constrained ML tree first**

**monophyly_test.pl**

likelihood of two topologies with given locus are significantly different, so given locus do not follow predefined group, discarded

likelihood difference of two topologies with given locus are not significant, so given locus follows predefined group, kept

# Filter for other purpose

## Filter out loci follow the molecular clock hypotheses

ATCGTAGGGCTGGCTAGTCGTAGCTA
ATCGTAGGGCTGGCGAGTCGT-GCTA
ATCGTAGGACTGGCTAGTCGTAGCTA
ATCGTACGGCTGGCTAGTCGTAGCTA
ATCGTAGGGCTGGCTAGTCGTAGCTA

**ML tree**

**Tree likelihood with molecular clock constrain**

**Tree likelihood without molecular clock constrain**

**Likelihood ratio test**

whether p > 0.05

**No**

**Significant difference when molecular clock hypotheses is applied, discarded**

**Yes**

**No significant difference between likelihood when applied with molecular clock hypotheses or not, kept**

**clocklikeness_test.pl**

# Detect cross-sample contamination

Genetic distance of diverged taxa in most
loci cannot be very close

detect_contamination.pl

# Detect cross-sample contamination

```
sample1 sample2 sample3 sample4 sample5 sample6
sample7 sample9
sample8 sample10 sample3
sample11 sample12
```

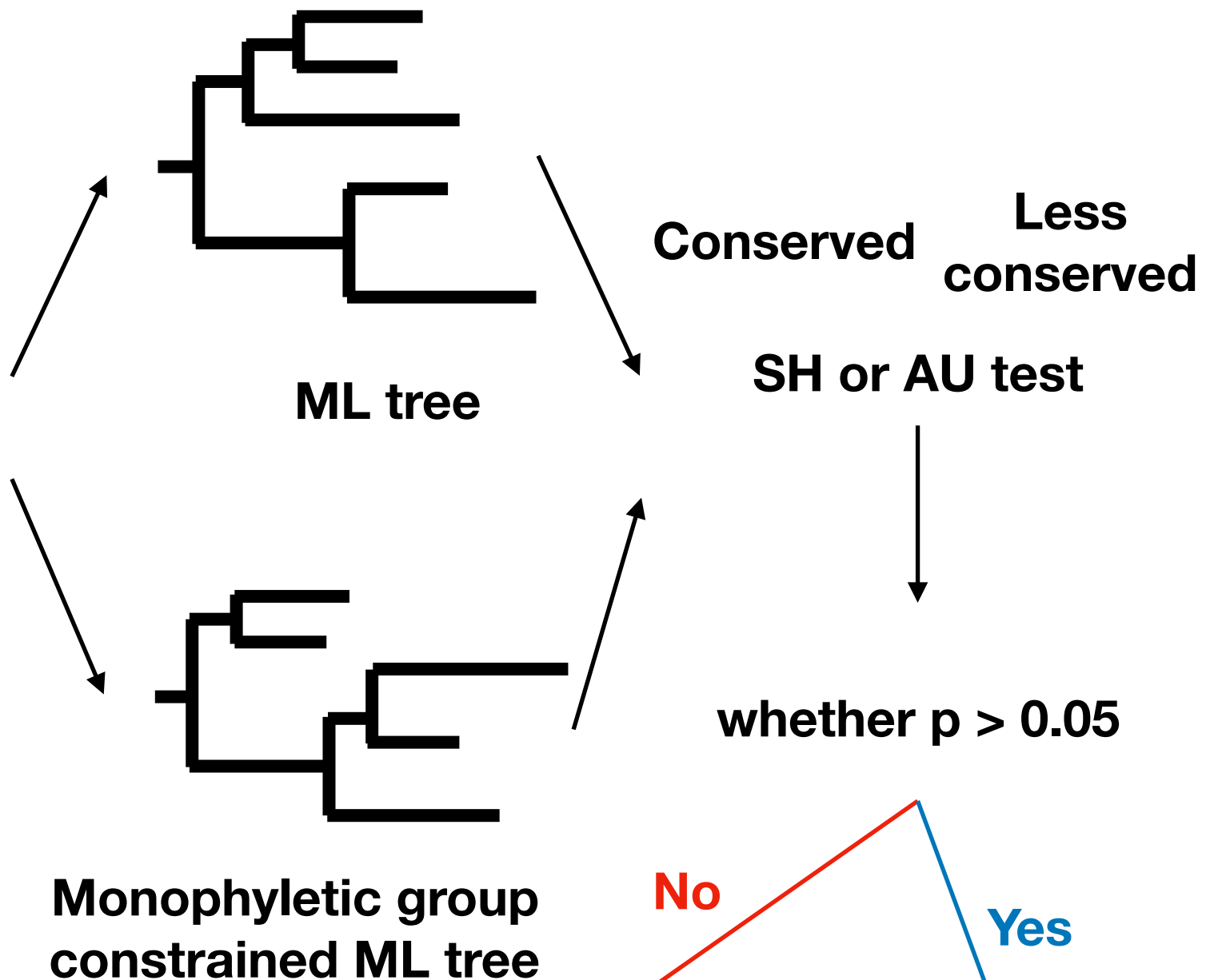**Define close related group in a txt file, one group a line**

**Permit <span style="color:red">one</span> sample in a group**

# Detect cross-sample contamination

**Closely related group1:  Human Chimp Orangutan**

↕

**too close p-distance between taxa (<= 0.002)**

**Closely related group2: Tilapia Zebra fish**

**Contamination rate (%) =**
**Potentially contaminated pair/Co-existence of this pair appeared in all loci*100**

**4434 loci in total**
**Potentially contamination between Human and Tilapia among all loci: 10 times**
**Co-existence of  Human-Tilapia pair among all loci: 2000 times**
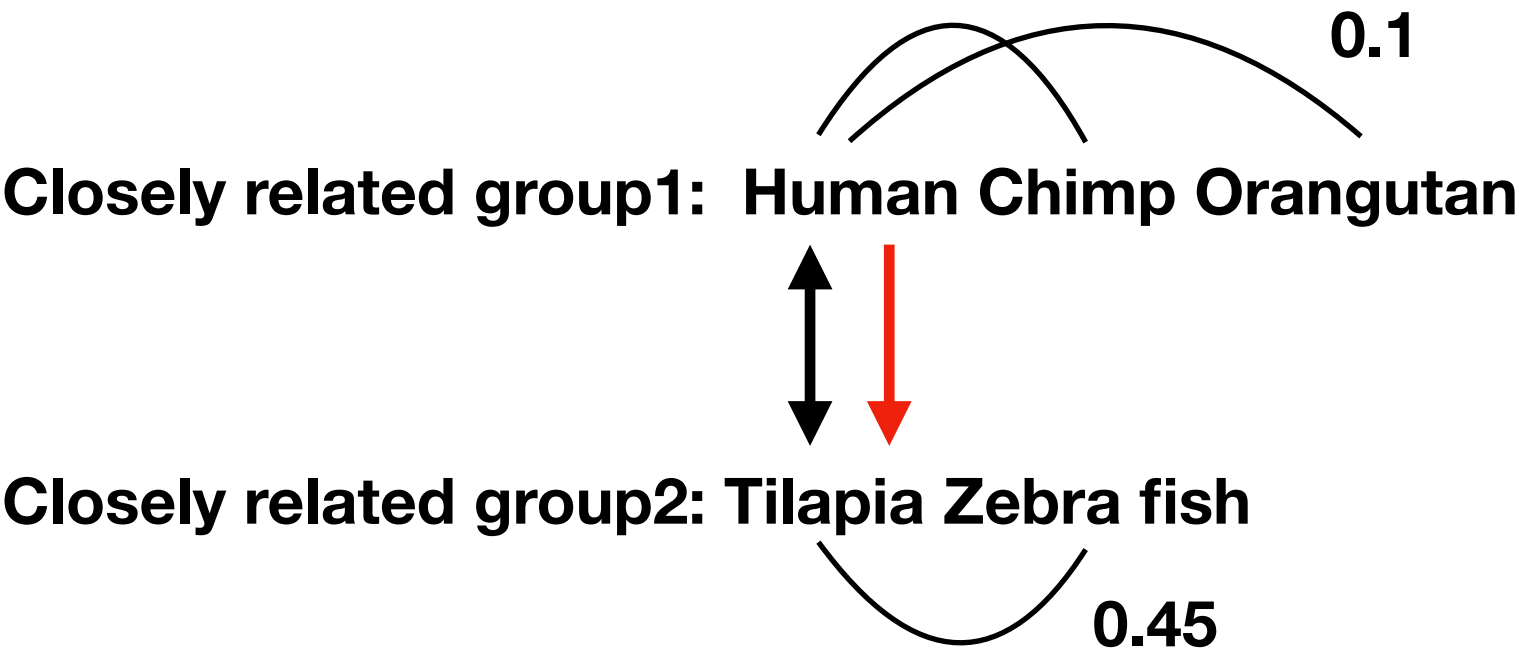
**Contamination rate of Human-Tilapia (%) = 10/2000*100 = 0.5%**

**Default threshold:**
**>=50% Contamination rate**
**>= 100 appearing time among all loci**

# Detect cross-sample contamination

0.1

Closely related group1:  Human Chimp Orangutan

Closely related group2: Tilapia Zebra fish

0.45

Human <- Tilapia: 10
Human -> Tilapia: 200

Sample of Human contaminated Tilapia

# Summary statistics

Summarized statistics for each **locus** including:
(1) Average length of coding region
(2) Average length of flanking region
(3) Length of alignment
(4) Average GC content
(5) Percentage of Missing data
(6) Pairwise distance

Summarized statistics for each **sample** including:
(1) Average length of captured sequences
(2) Average GC content
(3) Number of captured loci

**statistics.pl**

# Data filtering paradigm for coding region

**Contigs targeting the same loci**

**Add sequences from existing genomes if necessary** — get_orthologues.pl

**ILS?**
**Excessively trimmed?**
**Too few data?**

Filter data according to completeness level — pick_taxa.pl

**taxa>=80% & loci>=500**

**No** → remove loci <= 4 taxa

**Yes** → remove unqualified taxa

MSA — mafft_aln.pl

MSA — mafft_aln.pl

MSA — mafft_aln.pl

Gene tree — construct_tree.pl

Gene tree based Species tree reconstruction

ML Concatenated tree — concat_loci.pl

Gene tree based Species tree reconstruction — construct_tree.pl

**As expectation & high node support**

**Yes** → AM AWESOME ← **Yes**

**As expectation & Consistent**

**No** → Filter or trim loci

filter.pl, detect_contamination.pl, monophyly_test.pl, clocklikeness_test.pl

**No**

Filter or trim loci — filter.pl, detect_contamination.pl, monophyly_test.pl, clocklikeness_test.pl

No enough data for tree estimation

**Trouble shooting**

**More data required**

**Trouble shooting**

# Conventional analysis paradigm

**Inter**specific analysis

**Intra**specific analysis

Summarized statistics `statistics.pl`

Filtered coding sequences

Unaligned coding sequences with flanks

Concatenate all loci `concat_loci.pl`

Build gene trees `construct_tree.pl`

Aligned coding sequences with flanks `mafft_aln.pl`

Partition `Partition Finder`

Species tree `ASTRAL`

Filtered coding sequences with flanks `flank_filter.pl`

Concatenated tree `RAxML`

Generate consensus reference `consensus.pl`

Mapping statistics `map_statistics.pl`

Mapping `1-8 step of GATK.sh`

Trimmed reads

SNP calling `GATK.sh`

Reformat vcf file `vcftosnps.pl`

PCA   STRUCTURE   BEAST